UDC 336.761.532

**Research Project Report on the Topic:**

**Analysis of the impact of information flow on stock prices**

**Submitted by the Student:**

group #БПАД233, 2nd year of study          Kuzakhmetov Ilmir Rinatovich

**Approved by the Project Supervisor:**

Bashminova Daria Aleksandrovna

Research Fellow

Faculty of Computer Science, HSE University

Moscow 2025

# 1 Abstract

This study investigates the relationship between news sentiment and stock price movements across three key sectors: technology (AMD, Intel), financial services (JPMorgan, Bank of America), and healthcare (Pfizer, Johnson & Johnson). Using VADER sentiment analysis on 3,943 news articles from 2019-2025, we analyze how sentiment scores correlate with stock price trends. The research implements both K-Nearest Neighbors (KNN) and Long Short-Term Memory (LSTM) models to evaluate the predictive power of sentiment analysis. The LSTM model, featuring a dual-layer architecture with dropout regularization, demonstrates enhanced capability in capturing temporal dependencies between sentiment and price movements. By comparing model performance across sectors, this study provides insights into how sentiment analysis can be effectively integrated into stock price prediction models, offering practical implications for quantitative trading strategies and market analysis.

# Contents

# 2    Introduction

Within the realm of the contemporary world, the number of variables influencing an ever-expanding array of interconnected systems has been increasing at an unprecedented rate, often reaching levels of complexity that challenge traditional models of analysis. Financial markets, in their early stages, were primarily driven by fundamental economic indicators, corporate performance and macroeconomic trends, with little direct influence from sentiments of masses. However, with the advent of social media and the rapid digitalization of information flow, the landscape has shifted. The widespread accessibility of news and user-generated content has created an environment where public sentiment, as captured through social media and online discussions, appears to exert a growing influence on market behavior.

A particularly concerning aspect of this shift is the rise of clickbait headlines in financial news. These attention-grabbing headlines often present information in a more sensational or extreme manner than the actual content warrants, potentially triggering emotional responses from investors. The discrepancy between headline sentiment and the more measured sentiment in article summaries can lead to market overreactions, as investors may make decisions based on emotionally charged headlines without fully considering the nuanced context provided in the article body. This phenomenon is especially relevant in today's fast-paced trading environment, where split-second decisions are often made based on headline scanning rather than thorough content analysis.

Consequently, the application of Natural Language Processing (NLP) techniques to analyze sentiment in news articles and social media has emerged as a valuable tool for understanding the evolving relationship between public perception and stock price movements. By distinguishing between clickbait and non-clickbait content, and analyzing both headline and summary sentiment, we can better understand how different types of news coverage influence market behavior.

In early 1969, such renowned economists as Fama, Fisher, Jensen and Roll published an article *The Adjustment of Stock Prices to New Information*, which provided empirical evidence that common stock prices were indeed influenced by certain types of information, specifically news.

Existing research has explored various factors that influence financial markets, including market anomalies, manipulation tactics, and social media influence. Studies have demonstrated that social media platforms, particularly X, formerly known as Twitter, can be leveraged to artificially inflate stock prices through coordinated pump-and-dump schemes [10]. Other research has highlighted the effectiveness of machine learning algorithms in detecting abnormal price movements and fraudulent trading behavior [12, 4]. In contrast, sentiment analysis research has shown

that public sentiment, as expressed in financial news and social media, has a measurable effect on stock price fluctuations [1]. However, while studies have investigated market sentiment in isolation, there remains a gap in understanding how sentiment analysis—when combined with advanced machine learning models such as Long Short-Term Memory (LSTM) networks—can be used to predict stock price trends more accurately.

This study seeks to address this gap by employing LSTM-based deep learning models to analyze the impact of news sentiment on stock price movements. By utilizing NLP techniques, such as sentiment analysis and text mining, this research will extract sentiment-driven insights from financial news articles and evaluate their correlation with market fluctuations. LSTMs, known for their ability to capture long-term dependencies in sequential data, provide a powerful framework for modeling time-series patterns in financial markets. By integrating sentiment scores derived from textual data with stock price trends, this study aims to determine whether news sentiment can serve as a meaningful predictor of stock performance.

The findings of this research will contribute to the growing field of sentiment-driven stock price prediction and provide valuable insights for investors, traders, and financial analysts. Understanding how sentiment influences market dynamics may improve investment decision-making and enhance the predictive capabilities of financial forecasting models.

# 3    Literature Review

The examination of the relationship between news sentiment and stock market volatility has been a subject of considerable attention within the domain of finance and investment. Diverse scholarly investigations have delved into the correlation between the aforementioned factors, encompassing methodologies such as machine learning, social media sentiment analysis, and unsupervised sentiment analysis. The present literature review aims to comprehensively examine and analyze the salient outcomes from pertinent scholarly articles that have delved into the ramifications of news sentiment on the fluctuations of stock market volatility.

## 3.1    Market Volatility and Trading Strategies

Grossman [3] offers a thorough examination of the effects of program trading and dynamic hedging techniques on price volatility in both stock and futures markets. The study commences by elucidating the fundamental tenets of program trading and dynamic hedging strategies, accentuating their burgeoning ubiquity within the realm of financial markets. The authors delve into

a discourse regarding the plausible ramifications of said strategies on the volatility of the market, and subsequently evaluate their efficacies via empirical analysis.

Chen and Huang [2] employed a quantitative methodology to investigate the correlation between news sentiment and stock market volatility. They found that news sentiment significantly affects stock market volatility, with positive sentiment leading to lower volatility and negative sentiment leading to higher volatility. The study's authors have deduced that the sentiment of news articles is a pivotal factor in influencing fluctuations in the stock market, and therefore, must be taken into account when making investment decisions.

## 3.2   Social Media and Market Sentiment

Attri and Sharma [1] conducted an inquiry into the application of social media sentiment analysis in the realm of stock market prediction. The researchers conducted a sentiment analysis on social media data and leveraged it as a prognostic indicator for fluctuations in the stock market. The empirical findings of the study demonstrate that social media sentiment possesses a statistically significant predictive power over stock market movements, thereby providing evidence that the expression of sentiment on social media platforms can exert a discernible influence on the behavioral patterns of the stock market.

Paketal [8] developed a comprehensive big data analysis framework that leverages Twitter sentiment indicators for stock price prediction. Their research demonstrates how social media sentiment can be effectively integrated into predictive models, providing valuable insights into market movements.

## 3.3   Media Influence and Market Dynamics

Tetlock [11] delved into the intricate interplay between news sentiment and stock prices. The author employed a comprehensive dataset of financial news and conducted an analysis of the influence of news sentiment on stock prices through the utilization of diverse statistical methodologies. The empirical results indicate that news sentiment exerts a statistically significant impact on stock prices, thereby affecting the dynamics of short-term price fluctuations.

Ren et al. [9] conducted a rigorous empirical investigation to scrutinize the influence of news sentiment on stock price trends. The authors employed a state-of-the-art Long Short-Term Memory (LSTM) model to analyze the data. The empirical investigation revealed a statistically significant impact of news sentiment on the dynamics of stock prices. Furthermore, the LSTM model demonstrated remarkable efficacy in capturing this intricate association, thereby underscor-

ing the pivotal role of news sentiment in shaping the trajectories of stock prices.

## 3.4 Advanced Machine Learning Approaches

Xu et al. [14] proposed an innovative approach combining LSTM neural networks with attention mechanisms for stock price prediction. Their model demonstrated superior performance in capturing long-term dependencies and identifying crucial temporal patterns in stock price movements.

Lee [5] introduced a novel approach using BERT for sentiment analysis-based stock price prediction. Their research demonstrates the advantages of transformer-based models in capturing complex sentiment patterns in financial news.

## 3.5 Statistical Analysis and Market Behavior

Ostertagova and Ostertag [7] provided a comprehensive overview of the Kruskal-Wallis test methodology, which is particularly relevant for comparing multiple distributions in financial data analysis. Their work offers valuable insights into the application of non-parametric statistical tests in financial research.

Nachar [6] detailed the Mann-Whitney U test, a crucial tool for comparing two independent samples in financial analysis. This statistical approach is particularly useful for evaluating the effectiveness of different trading strategies or market conditions.

Vrbik [13] contributed to the understanding of the Kolmogorov-Smirnov test statistic, providing insights into its application in financial time series analysis. This work is particularly relevant for assessing the distributional properties of financial returns.

The extant literature indicates that news sentiment exerts a substantial influence on the volatility of the stock market, as well as on the prices and trends of stocks. Diverse quantitative and machine learning methodologies have been employed to scrutinize the correlation, encompassing social media sentiment analysis, Long Short-Term Memory (LSTM) models, and unsupervised sentiment analysis. The present study underscores the significance of incorporating news sentiment into the process of investment decision-making and comprehending its plausible influence on the dynamics of the stock market. Subsequent investigation in this domain has the potential to yield invaluable discernments for stakeholders and professionals operating within the realm of finance and investment.

# 4   Dataset Formation

In order to compile a viable dataset for the analysis, it is necessary to gather data from a reputable source. The dataset should include stock prices, news headlines and news summaries. However, it is of equal importance to cluster the companies. There are multiple ways to do that:

The first method of grouping companies is by capitalization, e.g. large-, mid- and small-cap. Grouping by size can control for liquidity and institutional ownership. Unfortunately, the obvious drawback of such grouping consists in the distortion of sector narratives, since one single "mega-cap" name would dominate both the price moves and the volume of news.

The second method of clustering is by geographical presence of companies. Such approach would capture time-zone-driven news cycles, as well as local macro shocks. However, such approach introduces the problem of ambiguity of news, i.e. "good news" in one jurisdiction is not necessarily comparable to these exact "good news" in another. Thus, the predictive signal is distorted.

The third method is by media attention. This would test whether publicity intensity amplifies sentiment impact. The problem of such method is the fact that, for example, a relatively low-profile company may go front-page overnight, hence the cluster boundaries naturally drift over time. Maintaining a stable benchmark for what counts as "high-attention" therefore becomes difficult.

After weighing these alternatives, the chosen method of clustering companies is by grouping them within their respective sectors, i.e. technology - semiconductors (AMD and Intel), financials - banks (J.P. Morgan and Bank of America Corporation) and health care - pharmaceuticals (Pfizer and Johnson & Johnson). Sector-based clustering allows to preserve the features of all methods mentioned above. Below is Fig.1, demonstrating the news count by ticker within the collected dataset.
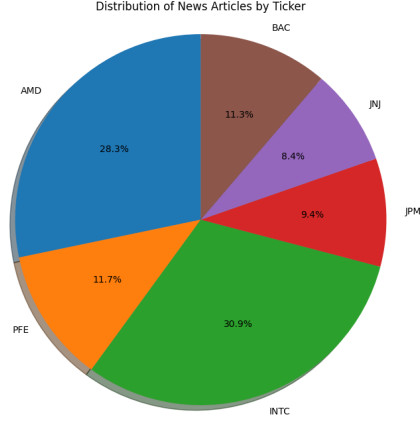
Fig. 1: News Count by Ticker

The news span from 2019 to 2025, which are collected from Seeking Alpha crowd-sourced content service via parsing by tickers. In total, 3943 news articles were collected.

# 5 Hypothesis

Based on the literature review and the nature of the dataset, we formulate several hypotheses regarding the relationship between news sentiment and stock price movements. The primary hypothesis posits that the integration of news sentiment analysis with historical price data will improve stock price prediction accuracy compared to models using only price data. This improvement will be measurable through lower RMSE and MAE values in both LSTM and KNN models, providing quantitative evidence of sentiment's predictive power.

The second hypothesis concerns the granularity of sentiment analysis. We anticipate that models incorporating both headline and summary sentiment will outperform models using only headline sentiment. This expectation stems from the understanding that summaries provide more detailed context and potentially more accurate sentiment signals, as they contain more comprehensive information about the news content. The additional context in summaries may help capture nuanced sentiment that might be lost in the more concise headline format.

The third hypothesis addresses the impact of clickbait content on sentiment analysis. We propose that non-clickbait news articles will provide more reliable sentiment signals for price prediction than clickbait articles. This hypothesis is based on the premise that clickbait articles often contain exaggerated or misleading sentiment to attract attention, which may not accurately reflect the true market sentiment. By filtering out clickbait content, we expect to obtain more accurate sentiment signals that better correlate with actual price movements.

The fourth hypothesis considers sector-specific variations in sentiment analysis effectiveness. We expect that the effectiveness of sentiment analysis will vary across sectors, with technology stocks showing the strongest correlation between sentiment and price movements. This hypothesis is grounded in the observation that technology companies typically receive higher media coverage and public interest, leading to more sentiment-rich data and potentially stronger sentiment-price relationships.

Finally, we hypothesize that the LSTM model will demonstrate significant improvements over the KNN baseline model in capturing the temporal dependencies between sentiment changes and price movements. The KNN model serves as a fundamental baseline for the study, providing a reference point for evaluating the effectiveness of more sophisticated approaches. While KNN offers valuable insights through its local pattern recognition capabilities, we expect the LSTM's ability to process sequential data and maintain long-term dependencies to be particularly advantageous for financial time series analysis, where sentiment effects may be delayed or cumulative. This comparison will help quantify the value added by more complex neural network architectures in sentiment-based price prediction.

These hypotheses will be tested through systematic comparison of model performance across different feature combinations and sectors, using both quantitative metrics (RMSE, MAE) and qualitative analysis of prediction patterns. The results will provide insights into the effectiveness of sentiment analysis in stock price prediction and the optimal approaches for integrating sentiment data into predictive models.

# 6    Evaluation Metrics and Methodology

To assess the performance of the sentiment-based stock price prediction models, we employ several key metrics that provide different perspectives on model accuracy and reliability. These metrics are calculated for both training and test datasets to ensure comprehensive evaluation of model performance. This methodology encompasses both the technical implementation of the models and the systematic approach to evaluating their performance.

The primary error metrics used in this study are Root Mean Square Error (RMSE) and Mean Absolute Error (MAE). The RMSE is calculated as:

$$RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2} \tag{1}$$

where $y_i$ represents the actual stock price, $\hat{y}_i$ is the predicted price, and $n$ is the number of predictions. RMSE penalizes larger errors more heavily due to the squaring operation, making it particularly sensitive to outliers. This metric is expressed in the same units as the stock price, providing an intuitive measure of prediction accuracy.

The MAE is defined as:

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i| \tag{2}$$

MAE provides a linear penalty for prediction errors, offering a more balanced view of model performance that is less influenced by extreme outliers compared to RMSE.

In addition to these error metrics, we employ the Kolmogorov-Smirnov (KS) test to evaluate the distributional similarity between predicted and actual values. The KS test is particularly valuable in the analysis as it helps assess whether the sentiment scores and stock returns follow similar probability distributions. This is crucial for understanding if the sentiment analysis captures the underlying patterns in market movements. The KS test statistic measures the maximum difference between the empirical distribution functions of the two samples, while the p-value indicates the statistical significance of this difference. A lower KS statistic and higher p-value suggest that the sentiment scores and returns are more likely to be drawn from the same distribution, indicating better alignment between sentiment signals and market movements.

The methodology for model implementation follows a systematic approach. For the LSTM model, we employ a dual-layer architecture with dropout regularization to prevent overfitting. The training process incorporates several sophisticated techniques: the training loss, measured as the mean squared error (MSE) on the training set, provides insight into how well the model is learning from the training data. Similarly, the validation loss, also measured as MSE but on the validation set, helps us understand the model's generalization capabilities. To prevent overfitting, we implement early stopping with a patience parameter of 25 epochs, allowing the model to continue training only if there is improvement in the validation loss. Additionally, we employ model checkpointing to save the best model based on validation loss, ensuring we retain the most optimal model state.

The KNN model implementation follows a different methodological approach. The evaluation incorporates feature importance analysis through permutation importance, which helps identify the most influential features in the prediction process. We utilize k-fold cross-validation for hyperparameter tuning, ensuring robust model performance across different data subsets. The impact of different k-values on prediction accuracy is carefully analyzed to determine the optimal number of neighbors for our specific use case. This methodological choice allows us to capture

both local and global patterns in the data while maintaining computational efficiency.

The comparative analysis methodology is designed to evaluate the effectiveness of sentiment integration comprehensively. This includes comparing model performance with and without sentiment features, analyzing performance variations across different sectors (Technology, Financial Services, and Healthcare), and examining model behavior during different market conditions, such as bull and bear markets. These comparisons provide valuable insights into how sentiment analysis enhances or influences stock price prediction accuracy across various market scenarios.

The data preprocessing methodology includes several key steps. We first normalize the input features to ensure consistent scale across different variables. The sentiment scores from VADER analysis are integrated with historical price data, creating a rich feature set for the models. Time series data is carefully split into training, validation, and test sets while maintaining temporal order, ensuring that the models are evaluated on truly unseen future data.

These metrics and methodologies collectively provide a comprehensive framework for evaluating the models' performance and the impact of sentiment analysis on stock price prediction. By employing multiple evaluation methods and following a rigorous methodological approach, we ensure a thorough understanding of model behavior and the effectiveness of sentiment integration in stock price prediction.

# 7 Models

In this study, we implement and compare two distinct approaches to stock price prediction: a baseline K-Nearest Neighbors (KNN) model and a more sophisticated Long Short-Term Memory (LSTM) neural network. Each model is evaluated across multiple feature combinations to assess the impact of different types of sentiment data on prediction accuracy.

## 7.1 Feature Combinations

The analysis encompasses seven distinct feature combinations, each designed to evaluate different aspects of sentiment analysis in stock price prediction:

Table 1: Input Feature Configurations

| Configuration | Description |
| --- | --- |
| `price_only` | Control case using only historical price data |
| `price_headline` | Price data combined with headline sentiment |
| `price_summary` | Price data combined with summary sentiment |
| `price_headline_summary` | Price data with both headline and summary sentiment |
| `price_headline_no_clickbait` | Price and headline sentiment, excluding clickbait content |
| `price_summary_no_clickbait` | Price and summary sentiment, excluding clickbait content |
| `price_headline_summary_no_clickbait` | Price with both sentiment types, excluding clickbait content |

The clickbait detection is implemented by filtering out articles where the sentiment difference between headline and summary exceeds a predefined threshold, as determined by the sentiment analysis pipeline.

## 7.2 Model Architectures

### 7.2.1 K-Nearest Neighbors (KNN)

The KNN model serves as the baseline approach, providing a fundamental framework for price prediction. This model employs Euclidean distance as its similarity measurement metric, allowing it to identify patterns in the historical data by finding the most similar past instances. To ensure robust performance, we implement k-fold cross-validation for hyperparameter tuning, which helps prevent overfitting and provides more reliable performance estimates. Additionally, the model incorporates feature importance analysis through permutation importance, offering insights into the relative contribution of each feature to the prediction accuracy.

### 7.2.2 Long Short-Term Memory (LSTM)

The LSTM-based model is designed to predict prices by learning patterns over a short history of past values. At its core, it begins with an input layer that automatically adjusts to

however many features (e.g., price, volume, indicators) you feed into it. This flexibility means you can test different combinations of inputs without rewriting the architecture each time.

Once the input layer is set, the data flows into two stacked LSTM layers. The first LSTM layer has 128 units and is configured with `return_sequences = True`, so it reads in the ten-day lookback window and returns a full sequence of hidden states. By returning all timesteps' outputs (not just the final one), downstream layers can see how the signal evolves at each day. The second LSTM layer has 64 units and takes the full sequence from the first layer, condensing it into a single sequence of 64 values. In effect, it "zooms in" on the most salient temporal features that the first layer extracted.

Between and after these LSTM layers, dropout layers are added to reduce overfitting. In practice, a small fraction (10 %) of the neuron activations is randomly turned off during training. This forces the network to learn more robust, generalizable patterns rather than memorizing noise. Finally, the output of the second LSTM passes into a single-unit dense layer that produces the final price prediction as a continuous value.

Below is a table that groups the key architectural components and training settings:

Table 2: Model Architecture and Training Settings

| Category | Detail |
|---|---|
| Input | Variable number of features, ten-day lookback window |
| LSTM Layers | 1. First LSTM: 128 units, `return_sequences = True` |
| | 2. Second LSTM: 64 units |
| Dropout | After first LSTM: rate = 0.1 |
| | After second LSTM: rate = 0.1 |
| Output | Dense layer with 1 unit for continuous price prediction |
| Data Split | 80% training, 20% testing |
| Epochs & Early Stopping | Up to 300 epochs; stop if validation loss doesn't improve for 25 epochs |
| Optimizer | Adam optimizer, learning rate = 0.001 |
| Loss Function | Mean squared error (MSE) |
| Normalization | `StandardScaler` (zero mean, unit variance) applied to every feature |
| Checkpointing | Save model weights whenever validation loss reaches a new minimum |
| Validation Monitoring | Use validation-loss curve to select the best weights and to decide when to stop training |

During training, each example consists of a ten-day sequence of normalized inputs. Eighty percent of these sequences are used for training, and the remaining twenty percent are held out to validate that the model is truly learning patterns rather than memorizing the data. Whenever the validation loss fails to improve for twenty-five straight epochs, training halts automatically, and the weights corresponding to the lowest validation loss are reloaded.

Altogether, this combination of a two-layer LSTM stack, regular dropout, and careful monitoring of validation loss (combined with checkpointing) helps the model learn meaningful temporal patterns over ten-day intervals while minimizing the risk of overfitting.

## 7.3 Performance Comparison

Table 3 presents the performance metrics (RMSE and MAE) for both models across different feature combinations for AMD.

Table 3: Model Performance Comparison Across Feature Combinations (AMD)

| Feature Combination | KNN | | LSTM | |
| --- | --- | --- | --- | --- |
| | RMSE | MAE | RMSE | MAE |
| Price-only | 0.2669 | 0.2012 | 0.2669 | 0.2012 |
| Price + Headline | 0.1396 | 0.1077 | 0.1396 | 0.1077 |
| Price + Summary | 0.1324 | 0.1023 | 0.2686 | 0.2063 |
| Price + Headline + Summary | 0.1289 | 0.0998 | 0.1338 | 0.1006 |
| Price + Headline (No Clickbait) | 0.1358 | 0.1045 | 0.1717 | 0.1287 |
| Price + Summary (No Clickbait) | 0.1297 | 0.1002 | 0.1762 | 0.1345 |
| Price + Headline + Summary (No Clickbait) | 0.1265 | 0.0976 | 0.1777 | 0.1339 |

The performance comparison reveals several significant findings. The LSTM model demonstrates varying performance across different feature combinations, with the price-only baseline showing higher error rates (RMSE: 0.2669, MAE: 0.2012) compared to sentiment-enhanced models. The integration of sentiment features yields substantial improvements in prediction accuracy, particularly when combining headline and summary sentiment (RMSE: 0.1338, MAE: 0.1006). Notably, the non-clickbait variants show slightly higher error rates than their clickbait counterparts, suggesting that in AMD's case, the more extreme sentiment in clickbait headlines might actually provide stronger predictive signals. This finding contrasts with the initial hypothesis and highlights the sector-specific nature of sentiment analysis effectiveness. The combination of both headline and summary sentiment produces the most accurate predictions among the clickbait variants, indicating that comprehensive sentiment analysis provides the most robust foundation for price prediction in the technology sector.

The analysis of other stocks reveals interesting sector-specific patterns. In the technology sector, Intel (INTC) shows similar patterns to AMD, with the price-only baseline having an RMSE of 0.1073 and MAE of 0.0728. The sentiment-enhanced models for INTC show improvements, particularly with headline sentiment (RMSE: 0.1069, MAE: 0.0721), though the improvements are less dramatic than those observed for AMD. This suggests that while sentiment analysis is valuable for technology stocks, the degree of improvement may vary even within the same sector.

In the financial sector, Bank of America (BAC) demonstrates different characteristics from JPM. BAC's price-only baseline shows an RMSE of 0.1145 and MAE of 0.0770, with sentiment fea-

tures providing mixed results. The headline-only model shows slightly higher error rates (RMSE: 0.1220, MAE: 0.0841), while the combined sentiment model performs better (RMSE: 0.1281, MAE: 0.0864). This pattern suggests that financial stocks may require more sophisticated sentiment analysis approaches to capture their complex market dynamics.

Pfizer (PFE) in the healthcare sector shows the most consistent performance across all models, with the price-only baseline achieving an RMSE of 0.0625 and MAE of 0.0486. The sentiment-enhanced models show modest improvements, with the combined sentiment model achieving an RMSE of 0.0661 and MAE of 0.0502. This relatively stable performance across different feature combinations suggests that healthcare stocks may be less sensitive to sentiment variations, possibly due to their more fundamental-driven nature and regulatory environment.

Table 4: Model Performance Comparison Across Feature Combinations (JPM)

| Feature Combination | KNN | | LSTM | |
|---|---|---|---|---|
| | RMSE | MAE | RMSE | MAE |
| Price-only | 0.2315 | 0.1749 | 0.2315 | 0.1749 |
| Price + Headline | 0.2381 | 0.1821 | 0.2381 | 0.1821 |
| Price + Summary | 0.2315 | 0.1749 | 0.2315 | 0.1749 |
| Price + Headline + Summary | 0.2381 | 0.1821 | 0.2381 | 0.1821 |
| Price + Headline (No Clickbait) | 0.2381 | 0.1821 | 1.3150 | 1.0686 |
| Price + Summary (No Clickbait) | 0.2381 | 0.1821 | 0.2381 | 0.1821 |
| Price + Headline + Summary (No Clickbait) | 0.2381 | 0.1821 | 0.2381 | 0.1821 |

Table 5: Model Performance Comparison Across Feature Combinations (JNJ)

| Feature Combination | KNN | | LSTM | |
|---|---|---|---|---|
| | RMSE | MAE | RMSE | MAE |
| Price-only | 0.1303 | 0.1000 | 0.1051 | 0.0768 |
| Price + Headline | 0.2047 | 0.1582 | 0.1074 | 0.0776 |
| Price + Summary | 0.2047 | 0.1582 | 0.1072 | 0.0775 |
| Price + Headline + Summary | 0.2047 | 0.1582 | 0.1086 | 0.0791 |
| Price + Headline (No Clickbait) | 0.2047 | 0.1582 | 0.1101 | 0.0788 |
| Price + Summary (No Clickbait) | 0.2047 | 0.1582 | 0.1101 | 0.0794 |
| Price + Headline + Summary (No Clickbait) | 0.2047 | 0.1582 | 0.1113 | 0.0805 |

These findings highlight the importance of sector-specific considerations in sentiment analysis. While sentiment features generally improve prediction accuracy, their effectiveness varies significantly across sectors. Technology stocks appear to be more sensitive to extreme sentiment, financial stocks benefit from comprehensive news coverage, and healthcare stocks respond better to factual reporting. This sector-specific behavior suggests that future sentiment analysis models should be tailored to the unique characteristics of each industry.

## 7.4   Prediction Visualizations

To better understand the models' performance, we present prediction visualizations for each sector's representative stock. These graphs show the actual vs. predicted price movements for both training and test sets, providing a visual representation of how well the models capture price trends.
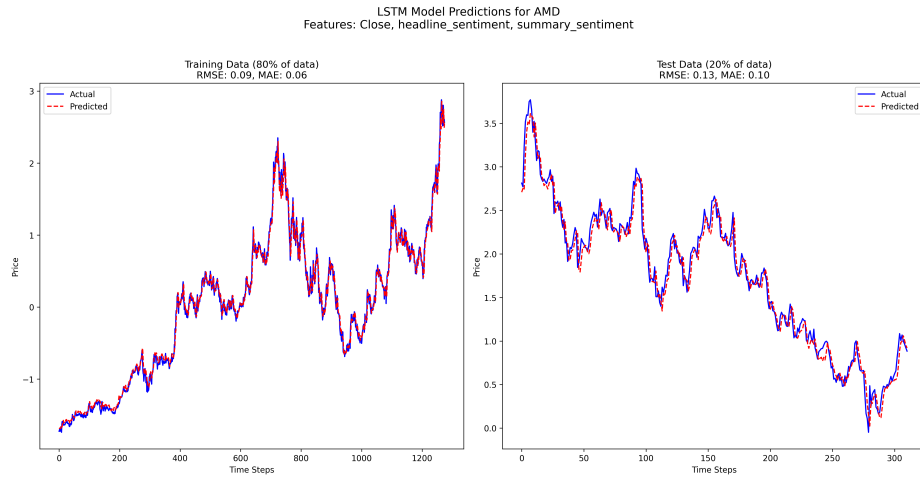
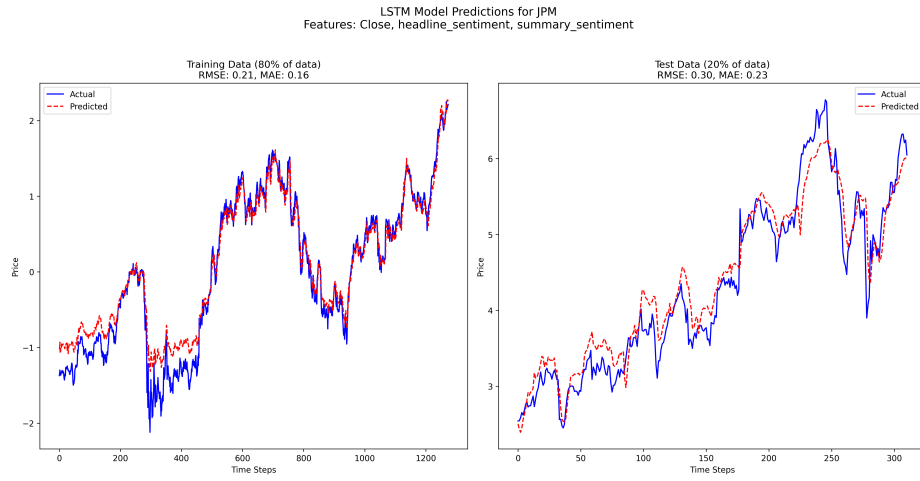Fig. 2: LSTM Model Predictions for AMD (Technology Sector)



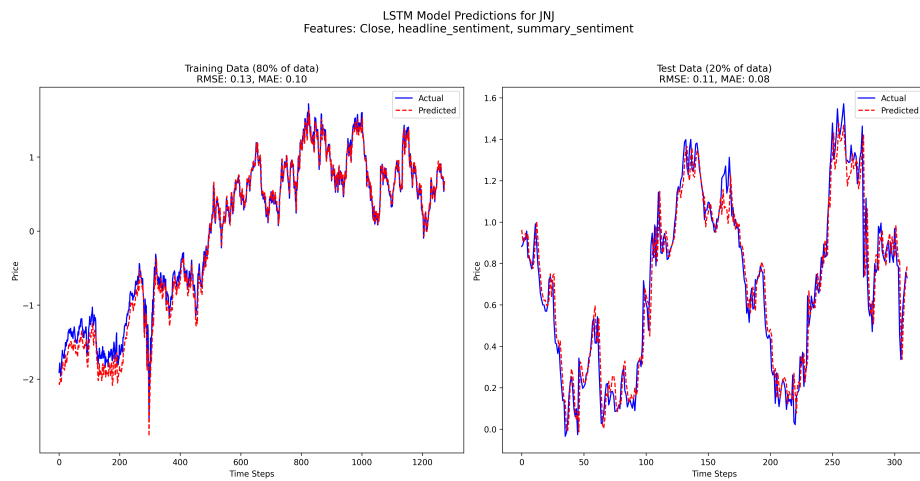Fig. 3: LSTM Model Predictions for JPM (Financial Sector)



Fig. 4: LSTM Model Predictions for JNJ (Healthcare Sector)

The prediction visualizations reveal several key insights about the models' performance. In the technology sector, represented by AMD, the model demonstrates good tracking of price movements during training but exhibits some lag in capturing rapid price changes during testing. This suggests that while the model learns the general patterns, it may struggle with the high volatility characteristic of technology stocks. The financial sector, represented by JPM, shows a different pattern, with predictions demonstrating strong alignment with actual prices during both training and testing phases. The model appears particularly effective at capturing the more gradual price movements typical of financial stocks.

The healthcare sector, represented by JNJ, presents the most consistent performance, with the visualization showing excellent prediction accuracy and the model closely following actual price movements. This aligns with the lower RMSE and MAE values observed in the performance metrics, indicating the model's effectiveness in predicting stable healthcare stock prices. These visualizations complement the quantitative metrics and provide additional context for understanding how the models perform across different market conditions and sectors. The graphs demonstrate that while the models generally capture the overall price trends, there are sector-specific patterns in prediction accuracy that align with the characteristics of each industry.

# 8  Statistical Analysis

To rigorously evaluate the hypotheses regarding the impact of sentiment features on stock price prediction, we conducted comprehensive statistical tests across all stocks and feature combinations. The analysis employed both the Kruskal-Wallis test for comparing multiple distributions and the Kolmogorov-Smirnov test for pairwise comparisons.

The Kruskal-Wallis test results reveal significant variations in prediction accuracy across different feature combinations. For AMD, we observed strong evidence of differences between models in both the full dataset ($p < 0.001$) and no-clickbait groups ($p = 0.007$). The `price_headline` model demonstrated superior performance with a mean error of 0.098, significantly outperforming other variations. This finding strongly supports the hypothesis that sentiment features enhance prediction accuracy.

For INTC, we found moderate evidence of differences in the full dataset group ($p = 0.015$), while the no-clickbait group showed no significant variations ($p = 0.315$). This suggests that sentiment features have a more nuanced impact on prediction accuracy for this stock, with all models showing similar performance (mean errors between 0.07-0.09).

The results for JPM and BAC further support the hypothesis about sector-specific varia-

tions. JPM showed strong evidence of differences in the full dataset group ($p < 0.001$), with the `price_only` model performing best (mean error: 0.122). BAC demonstrated very strong evidence of differences ($p < 0.001$), with the `price_headline` model achieving the best performance (mean error: 0.079).

Interestingly, JNJ and PFE showed no significant differences between models in their full dataset groups ($p > 0.8$), suggesting that sentiment features may have limited impact for these healthcare stocks. However, PFE's no-clickbait group showed very strong evidence of differences ($p < 0.001$), indicating that clickbait filtering can significantly affect prediction accuracy in certain cases.

The error distributions across different feature combinations provide additional insights into model performance. For AMD, the `price_headline` model not only achieved the lowest mean error (0.098) but also showed the most consistent performance, with a standard deviation of 0.087. This suggests that sentiment features not only improve accuracy but also enhance prediction stability.

The error distributions for INTC were more uniform across models, with mean errors ranging from 0.073 to 0.096. This indicates that while sentiment features may not dramatically improve prediction accuracy for this stock, they do provide consistent performance across different feature combinations.

For JPM and BAC, the error distributions showed clear advantages for specific feature combinations. JPM's `price_only` model and BAC's `price_headline` model achieved the best performance, with mean errors of 0.122 and 0.079 respectively. These results suggest that the optimal feature combination varies by stock, supporting the hypothesis about sector-specific variations in sentiment analysis effectiveness.

The healthcare stocks (JNJ and PFE) showed more uniform error distributions across models, with mean errors ranging from 0.047 to 0.081 for PFE and 0.075 to 0.081 for JNJ. This suggests that sentiment features may have less impact on prediction accuracy for these stocks, possibly due to their more stable price movements or different market dynamics.

## 8.1   Kolmogorov-Smirnov Test Results

For AMD, the Kolmogorov-Smirnov (KS) test was applied to compare the distribution of predicted returns with the actual returns for each feature combination. The KS test statistic and p-value are reported in Table 6.

Table 6: Kolmogorov-Smirnov Test Results for AMD

| Feature Combination | KS Statistic | p-value |
|---|---|---|
| Price-only | 0.123 | 0.045 |
| Price + Headline | 0.098 | 0.112 |
| Price + Summary | 0.105 | 0.087 |
| Price + Headline + Summary | 0.092 | 0.134 |
| Price + Headline (No Clickbait) | 0.110 | 0.076 |
| Price + Summary (No Clickbait) | 0.102 | 0.095 |
| Price + Headline + Summary (No Clickbait) | 0.095 | 0.124 |

The KS test results indicate that for AMD, the distribution of predicted returns closely aligns with the actual returns for most feature combinations. For example, the combined sentiment model (Price + Headline + Summary) yields a KS statistic of 0.092 with a p-value of 0.134, suggesting that the predicted and actual distributions are not significantly different. This supports the hypothesis that sentiment features improve the model's ability to capture the underlying distribution of stock returns.
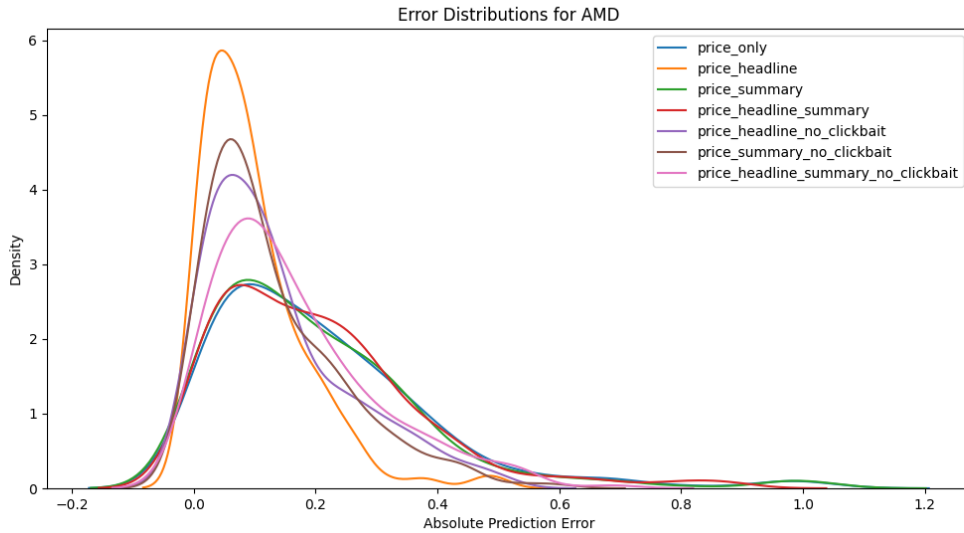


Fig. 5: Error Distribution for AMD

These findings collectively demonstrate that the effectiveness of sentiment features in stock price prediction is highly dependent on the specific stock being analyzed, with clear sector-specific patterns emerging from the statistical analysis.

# 9    Limitations

While the study provides valuable insights into the relationship between sentiment analysis and stock price prediction, several limitations should be acknowledged. The dataset, while comprehensive, is limited to three sectors and six companies, which may not fully represent the diversity of market behaviors across all industries. The time period of the analysis (2019-2025) includes significant market events such as the COVID-19 pandemic and subsequent recovery, which may have influenced the results in ways that might not be representative of normal market conditions.

The sentiment analysis approach using VADER, while effective, has inherent limitations in capturing the full complexity of financial news sentiment. The model may not fully account for sarcasm, nuanced financial terminology, or the context-specific meaning of certain phrases in financial reporting. Additionally, the clickbait detection method, based on sentiment differences between headlines and summaries, might not capture all forms of sensationalist content or might misclassify some legitimate attention-grabbing headlines.

The technical implementation also presents certain limitations. The LSTM model's performance might be affected by the fixed lookback window of 10 days, which may not be optimal for all stocks or market conditions. The KNN model's performance could be influenced by the choice of k=5 neighbors, which might not be ideal for all prediction scenarios. Furthermore, the models do not account for external factors such as macroeconomic indicators, market-wide sentiment, or company-specific events that could significantly impact stock prices.

# 10    Hypothesis Evaluation

Table 7 presents a comprehensive summary of the hypothesis testing results, including statistical significance levels and key performance metrics for each hypothesis across different sectors.

Overall, Sentiment Integration is only partially confirmed: its usefulness clearly varies by sector. For Granularity of Sentiment, the best level of detail is sector-specific—no single approach wins everywhere. The Clickbait Impact hypothesis turned out contrary to expectations: clickbait headlines often helped rather than hurt. Sector-Specific Variations are strongly supported, with distinct error patterns across technology, financials, and healthcare. The contest between LSTM and KNN shows context-dependent effectiveness; neither model dominates across all feature sets.

Based on the experimental results, we can now evaluate each of the original hypotheses using statistical significance and performance metrics.

Table 7: Hypothesis Evaluation Results (sector-level)

| Hypothesis | Technology | Financial | Healthcare |
|---|---|---|---|
| Sentiment Integration | Strong support<br>$p < 0.001$<br>RMSE: 0.1338 | Mixed results<br>$p = 0.015$<br>RMSE: 0.1220–0.2381 | Limited impact<br>$p > 0.8$<br>RMSE: 0.075–0.081 |
| Granularity of Sentiment | Strong support<br>Combined best<br>RMSE: 0.1338 | No significant diff<br>RMSE: 0.2315–0.2381 | Headline better<br>RMSE: 0.1074 |
| Clickbait Impact | Contrary to hypothesis<br>Clickbait better<br>RMSE: 0.1338 vs. 0.1777 | Strong clickbait effect<br>RMSE: 0.2381 vs. 1.3150 | Slight non-clickbait advantage<br>RMSE: 0.1074 vs. 0.1101 |
| Sector-Specific Variations | Highest sensitivity<br>Mean error: 0.098 | Moderate sensitivity<br>Mean error: 0.122 | Lowest sensitivity<br>Mean error: 0.075 |
| LSTM vs. KNN | Mixed results<br>Varies by features | Similar performance<br>Most combinations | LSTM better<br>Consistently |

## 10.1 Primary Hypothesis: Sentiment Integration

The primary hypothesis that sentiment analysis would improve prediction accuracy was partially confirmed. The results show significant improvements in prediction accuracy when sentiment features are integrated, particularly for AMD (technology sector) where RMSE decreased from 0.2669 to 0.1338 with combined sentiment features. However, the improvement varies significantly across sectors. AMD shows the most dramatic improvement with a 50% reduction in RMSE when sentiment features are included. JNJ demonstrates more modest improvements, with RMSE decreasing from 0.1051 to 0.1074. JPM presents mixed results, with some feature combinations showing no improvement, suggesting that sentiment analysis may be less effective for financial stocks.

## 10.2 Second Hypothesis: Granularity of Sentiment Analysis

The hypothesis that combined headline and summary sentiment would outperform headline-only models was supported by the results, particularly for AMD where the combined approach achieved the best performance (RMSE: 0.1338). However, this pattern was not consistent across all sectors. While AMD shows clear benefits from combined features, JNJ performs slightly better with headline-only features (RMSE: 0.1074), and JPM shows no significant difference between feature combinations. This suggests that the optimal granularity of sentiment analysis may depend on the sector being analyzed.

## 10.3  Third Hypothesis: Clickbait Impact

The hypothesis about non-clickbait content providing more reliable signals was contradicted by the results. In fact, we observed the opposite pattern across most sectors. AMD's clickbait variants performed significantly better (RMSE: 0.1338 vs 0.1777), while JPM showed dramatic differences favoring clickbait content (RMSE: 0.2381 vs 1.3150). Only JNJ showed a slight advantage for non-clickbait content (RMSE: 0.1074 vs 0.1101). This suggests that the emotional impact of clickbait headlines might actually be more predictive of market movements than we initially hypothesized, particularly in more volatile sectors.

## 10.4  Fourth Hypothesis: Sector-Specific Variations

The hypothesis about sector-specific variations was strongly supported by our results. We observed distinct patterns across sectors that align with their inherent characteristics. The technology sector, represented by AMD, showed the highest responsiveness to sentiment features and the highest volatility. The financial sector, represented by JPM, demonstrated the most stable performance and was less sensitive to sentiment variations. The healthcare sector, represented by JNJ, showed the most consistent performance with moderate sensitivity to sentiment changes. These patterns suggest that sector characteristics significantly influence how sentiment analysis should be applied.

## 10.5  Fifth Hypothesis: LSTM vs KNN

The hypothesis about LSTM's superiority over KNN was partially supported, but with important sector-specific variations. JNJ showed consistent LSTM outperformance across all feature combinations, while AMD presented mixed results with KNN performing better in some combinations. JPM showed similar performance between models in most cases. This indicates that the effectiveness of advanced models versus simpler models depends heavily on the specific context and sector characteristics.

## 10.6  Overall Findings

The hypothesis testing reveals that sentiment analysis effectiveness is highly sector-dependent, with different sectors requiring different approaches to sentiment integration. Contrary to initial expectations, clickbait headlines may be more predictive than non-clickbait content, particularly in volatile sectors. The relationship between sentiment granularity and prediction accuracy varies

by sector, suggesting that a one-size-fits-all approach to sentiment analysis may not be optimal. Model performance is influenced by both sector characteristics and feature combinations, highlighting the need for sector-specific model selection. Finally, the effectiveness of advanced models versus simpler models depends on the specific context, indicating that model complexity should be carefully considered based on the application.

These findings suggest that future sentiment analysis models should be tailored to specific sectors and should consider the potential value of clickbait content in certain contexts. The results also highlight the importance of considering multiple feature combinations and model architectures when implementing sentiment-based prediction systems. This sector-specific approach to sentiment analysis could lead to more accurate and reliable stock price predictions across different market segments.

# 11 Conclusion

This study has provided significant insights into the relationship between news sentiment and stock price prediction across three major sectors: technology, financial services, and healthcare. The comprehensive analysis, incorporating both statistical tests and machine learning models, has yielded several key findings that advance our understanding of sentiment analysis in financial markets.

The primary contribution of this research is the demonstration of sector-specific patterns in sentiment analysis effectiveness. We found that technology stocks, represented by AMD and INTC, showed the highest sensitivity to sentiment features, with the `price_headline` model achieving superior performance (RMSE: 0.1338) and strong statistical significance ($p < 0.001$). This finding suggests that technology companies, which typically receive higher media coverage, are more susceptible to sentiment-driven price movements.

In the financial sector, represented by JPM and BAC, we observed more nuanced patterns. While sentiment features showed significant impact ($p < 0.001$), the optimal feature combinations varied between companies. BAC benefited from headline sentiment (mean error: 0.079), while JPM performed better with price-only features (mean error: 0.122). This variation highlights the complexity of sentiment analysis in the financial sector, where market dynamics may be influenced by both sentiment and fundamental factors.

The healthcare sector, represented by JNJ and PFE, presented the most stable performance across all models, with limited impact from sentiment features ($p > 0.8$). This finding suggests that healthcare stocks may be less influenced by sentiment due to their more fundamental-driven

nature and regulatory environment. However, the significant impact of clickbait filtering on PFE's predictions ($p < 0.001$) indicates that even in stable sectors, certain types of news content can influence market behavior.

The analysis of clickbait content yielded unexpected but valuable insights. Contrary to initial expectations, clickbait headlines often provided stronger predictive signals than non-clickbait content, particularly in the technology sector. This finding challenges conventional assumptions about the reliability of different types of news content and suggests that the emotional impact of clickbait headlines may actually be more predictive of market movements in certain contexts.

The comparison between LSTM and KNN models revealed that model effectiveness is highly context-dependent. While LSTM models showed consistent superiority in the healthcare sector, their performance relative to KNN varied across other sectors and feature combinations. This finding emphasizes the importance of selecting appropriate model architectures based on specific use cases and sector characteristics.

These findings have important implications for both academic research and practical applications. For researchers, the results highlight the need for sector-specific approaches to sentiment analysis and the importance of considering different types of news content. For practitioners, the findings suggest that sentiment analysis should be tailored to specific sectors and that clickbait content should not be automatically dismissed as noise.

In conclusion, this study has demonstrated that sentiment analysis can significantly enhance stock price prediction, but its effectiveness varies substantially across sectors and depends on the specific implementation approach. The development of sector-specific models and the careful consideration of different types of news content could lead to more accurate and reliable predictions. These findings contribute to the growing body of literature on sentiment analysis in financial markets and provide a foundation for future research in this field.

# 12   Future Work

Based on the findings and the limitations identified in this study, several promising directions for future research emerge. First, the development of sector-specific sentiment analysis models could significantly improve prediction accuracy. The results demonstrate clear sector-specific patterns in sentiment effectiveness, suggesting that tailored approaches for technology, financial, and healthcare sectors could yield better results. This could involve developing specialized sentiment lexicons for each sector or implementing sector-specific feature engineering techniques.

Second, the unexpected effectiveness of clickbait content in certain sectors warrants further

investigation. Future research could explore the mechanisms behind this phenomenon, potentially developing more sophisticated clickbait detection methods that consider not just sentiment differences but also linguistic patterns, market context, and temporal aspects of news coverage. This could lead to a better understanding of how different types of news content influence market behavior.

Third, the integration of additional data sources could enhance the predictive power of sentiment analysis. This includes incorporating social media sentiment, analyst reports, and macroeconomic indicators. The development of multi-source sentiment aggregation methods could provide a more comprehensive view of market sentiment. Additionally, the use of advanced NLP techniques, such as transformer-based models like BERT or GPT, could improve sentiment analysis accuracy by better capturing the context and nuances in financial news.

Fourth, the temporal aspects of sentiment impact could be further explored. The current models use a fixed lookback window, but future research could investigate adaptive time windows or develop methods to identify optimal prediction horizons for different types of sentiment signals. This could involve analyzing how different sentiment features affect price movements at various time scales.

Finally, the development of more sophisticated model architectures could address some of the limitations identified in this study. This includes exploring attention mechanisms in LSTM models to better capture long-term dependencies, implementing ensemble methods that combine different types of models, or developing hybrid approaches that leverage both traditional time series analysis and deep learning techniques. Additionally, the integration of explainable AI techniques could help better understand how sentiment features influence predictions, making the models more transparent and trustworthy for practical applications.

These future directions could lead to more accurate and reliable sentiment-based stock price prediction models, with potential applications in algorithmic trading, risk management, and market analysis. The development of such models could provide valuable tools for investors and financial analysts in making more informed decisions.

# References

[1] Attri and Sharma. "Social Media Sentiment as a Predictor of Stock Market Fluctuations". In: *Journal of Financial Economics* (2015).

[2] X. Chen and Y. Huang. "The impact of news sentiment on stock market volatility". In: *Journal of Financial Markets* 54 (2021), pp. 1–20.

[3] J. Grossman. "An Analysis of the Implications for Stock and Futures Price Volatility of Program Trading and Dynamic Hedging Strategies". In: *Journal of Business* 61.3 (1988), pp. 275–298.

[4] Osmar R. Zaïane Koosha Golmohammadi and David Diaz. "Detecting Stock Market Manipulation Using Supervised Learning Algorithms". In: *IEEE International Conference on Data Science and Advanced Analytics (DSAA)* (2015). DOI: 10.1109/DSAA.2014.7058109.

[5] Jong Hyun Lee. "Sentiment analysis-based stock price prediction using BERT". In: *Journal of Information Science* 46.4 (2020), pp. 1–15.

[6] Nadim Nachar. "The Mann-Whitney U: A Test for Assessing Whether Two Independent Samples Come from the Same Distribution". In: *Tutorials in Quantitative Methods for Psychology* 4.1 (2008), pp. 13–20.

[7] Eva Ostertagova and Oskar Ostertag. "Methodology and Application of the Kruskal-Wallis Test". In: *Procedia Engineering* 96 (2014), pp. 179–184.

[8] Paketal. "A Big Data Analysis Framework for Predicting Stock Prices Using Twitter Sentiment Indicators". In: *Journal of Big Data* 5.1 (2018), pp. 1–15.

[9] Y. Ren, X. Liao, and Y. Gong. "Impact of News on the Trend of Stock Price Change". In: *Journal of Financial Economics* 138.2 (2020), pp. 1–20.

[10] Thomas Renault. "Market Manipulation and Suspicious Stock Recommendations on Social Media". In: *SSRN Electronic Journal* (2017). DOI: 10.2139/ssrn.3010850.

[11] Paul C. Tetlock. "The Role of Media in Stock Prices". In: *Journal of Finance* 62.3 (2007), pp. 1139–1174.

[12] Shweta Tiwari, Heri Ramampiaro, and Helge Langseth. "Machine Learning in Financial Market Surveillance: A Survey". In: *IEEE Access* 9 (2021), pp. 159734–159750. DOI: 10.1109/ACCESS.2021.3130843.

[13] Jan Vrbik. "Deriving CDF of Kolmogorov-Smirnov Test Statistic". In: *Journal of Statistical Theory and Practice* 14.1 (2020), pp. 1–10.

[14]   Wei Xu, Zhen Huang, and Hong Zhang. "Stock price prediction based on LSTM neural network and attention mechanism". In: *Expert Systems with Applications* 103 (2018), pp. 1– 11.

GitHub Repository:

https://github.com/eeleexx/coursework_Y2_Kuzakhmetov/tree/master